# *SESAME*

## (Video **SE**arch with **S**peed and **A**ccuracy for **M**ultimedia **E**vents)
## Multimedia Event Detection (MED) System

*November 26, 2012*

SRI International

UNIVERSITY OF AMSTERDAM
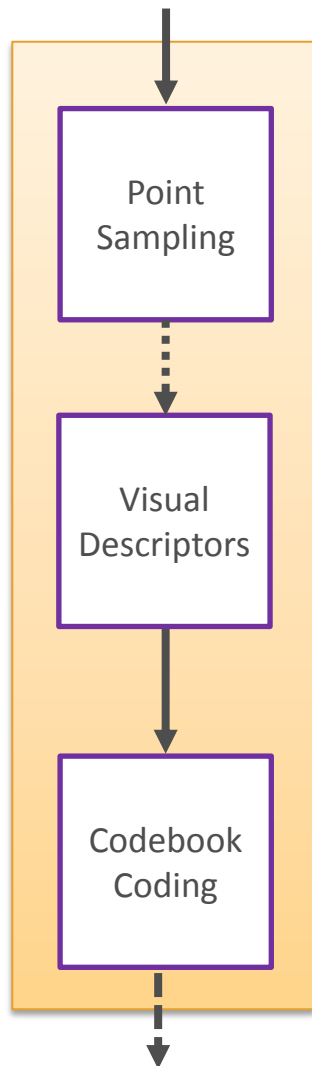
UNIVERSITY OF SOUTHERN CALIFORNIA · 1880

# Overview

- **Event classifiers**

- **Fusion and threshold selection**

- **Waypoint experiments on development set**

- **MED12 evaluation results**

# 14 Low-level and High-level Event Classifiers

- **Low-level features**
  - visual features (2)
  - motion features (5)
  - audio features (1)

- **Concept-level features:**
  - visual concepts (2)
  - ASR (2)
  - video OCR (2)

# Visual Features: Bag-of-Words and Difference Coding

Point Sampling

Visual Descriptors

Codebook Coding

Bag of words for
event agents *and*
visual scenes, objects, persons, actions

State-of-the-art

ColorSIFT [Van de Sande et al. TPAMI 2010]

Soft-Assignment [Van Gemert et al. TPAMI 2010]

Real-time Bag-of-Words [CIVR09 best paper]

TSIFT [under review]

# Two event classifiers based on visual features

1 frame sampled in every 2 seconds of video

| | Sampling | Descriptors | Codebook | Aggregation | Kernel |
|---|---|---|---|---|---|
| 1. Color average coding | Dense Harris | PCA reduced: SIFT, CSIFT, TSIFT | 4096, hard 1x1,1x3 | Average | Fast HIK |
| 2. Color difference soft coding | Dense | PCA reduced: SIFT, CSIFT, TSIFT | 1024, soft 1x1,1x3 | Average | Linear |

Waypoint experiments showed:
- average coding outperformed difference coding
- difference coding complemented average coding in late fusion experiments
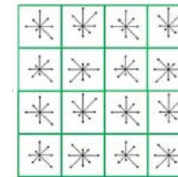
# Low-level motion features

- **STIP:**
  - Corner like detectors in 3D
  - 72-dim HOG + 90-dim HOF



- **MoSIFT**
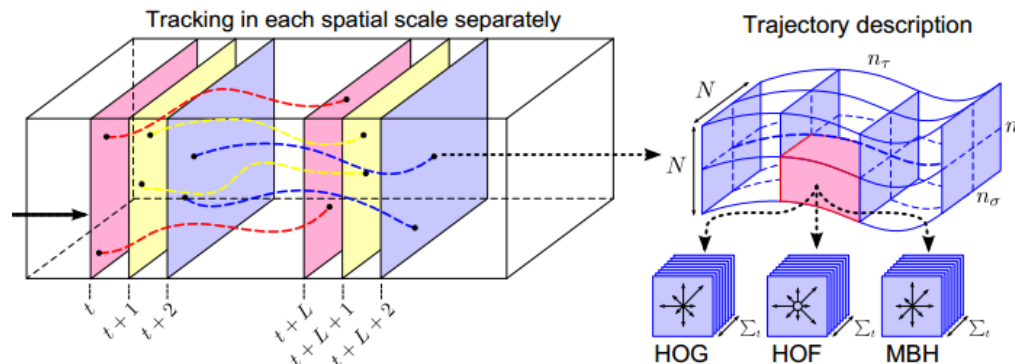  - SIFT like detectors in 2D, filtered by motion
  - Extracted, quantized and pooled by CMU



Alexander Kläser et al,
BMVC 2008

- **Dense Trajectories (DT):**
  - Generate tracklets for densely sampled points
  - Describe each tracklet by shape, HoG, HoF and MBH of the volume around it



Tracking in each spatial scale separately

Trajectory description

HOG   HOF   MBH

Wang et al.,
CVPR 2011

# Event classifiers using low-level motion features

- **5 event classifiers:**

| Event Classifier | Feature | Descriptor | Aggregation | Kernel |
|:---:|:---:|:---:|:---:|:---:|
| 1 | STIP | 1st-order Fisher | Average | Gaussian |
| 2 | STIP | 2nd-order Fisher | Average | Gaussian |
| 3 | DT | 1st-order Fisher | Average | Gaussian |
| 4 | DT | 2nd-order Fisher | Average | Gaussian |
| 5 | MoSIFT | MoSift | Average | χ2 |

- **Waypoint experiment showed:**
  - Dense Trajectory gives the best performance
  - 2nd order Fisher vector is better than 1st order
  - All 3 motion features are complementary in late fusion experiments

# Event Classifier using Low-level Audio Features

• Codebook size = 1000

```
                              ┌──────────────┐
                              │  Codebook    │
                              │  generation  │
                              └──────┬───────┘
                                     │
                                     ▼
┌──────────┐   ┌──────────────┐  ┌──────────────┐   ┌──────────────┐   ┌────────────┐
│  Video   │──▶│   Feature    │─▶│    Vector    │──▶│ Bag of audio │──▶│    SVM     │
│  files   │   │  extraction  │  │ quantization │   │    words     │   │ Classifier │
└──────────┘   └──────────────┘  └──────────────┘   └──────────────┘   └────────────┘
```

• No histogram normalization

• 16 kHz sampling rate

• MFCCs every 10 ms

• 12 coeff.+ log-energy + Δ+ Δ-Δ of each = 39 dim total

• No MFCC normalization

• Soft quantization
→ add distance from nearest codeword instead of +1 to the histogram for each quantized vector

• Histogram intersection kernel

# Event Classifiers using Visual Concept Detectors

- **1346 concept detectors**
  - 346 concepts from the TRECVID 2012 SIN task
  - 1,000 concepts from ImageNet
  - All trained using color difference coding with linear SVM

- **Two event classifiers**
  - One used random forests
  - One used a non-linear SVM

# Automatic Speech Recognition (ASR)

**audio** → Acoustic Feature Extraction → Supervised Speech/Nonspeech Segmentation (e.g., HMM-based) → ASR

*audio*

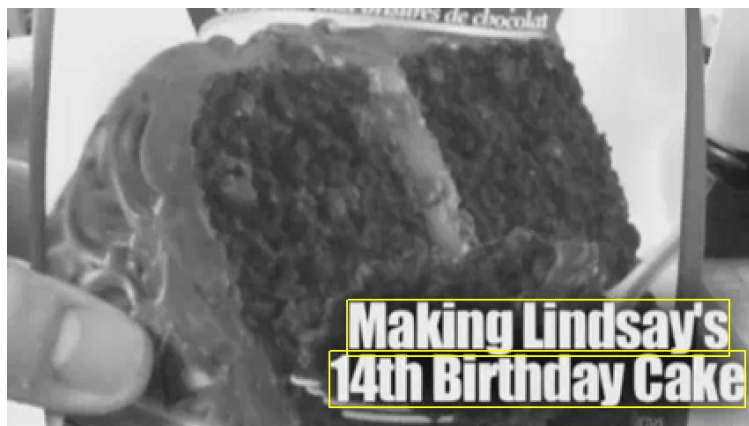*front-end features (frame-level)*

*"Meetings" ASR system*



- **Un-adapted ASR system trained on far-field microphone meetings data**

- **3-state ergodic HMM for audio segmentation (speech, music, other)**

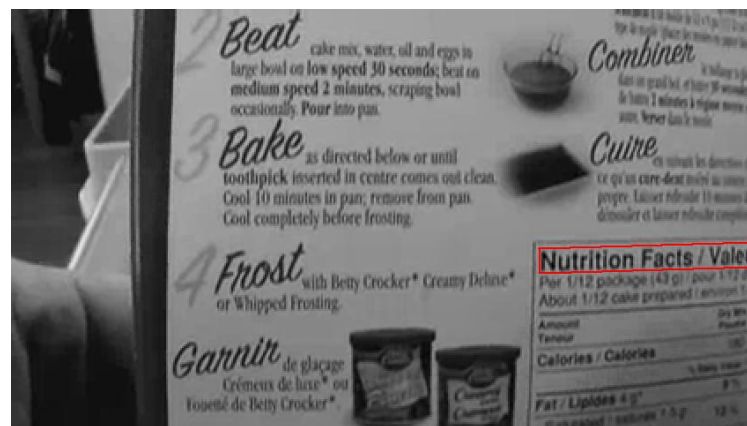- **ASR configured to recognize spoken English**

# Video OCR

- **SRI's video optical character recognition (video OCR) for detection, tracking, and recognition of text**
  - **recognizes both overlay text and in-scene text**
  - **configured to recognize English language text**

Text captions

In-scene text



**"Making Lindsay s 14th Birthday Cake"**

**"Nutrition Facts – Valeu"**

# 4 Event Classifiers for ASR and OCR Text

- **Each classifier measures the overlap of text in the test video with text in the event model using logistic regression**

- **Two event classifiers (one for ASR and one for OCR) based on text found in training set clips**
  - Unigram bag-of-words event models

- **Two event classifiers (one for ASR and one for OCR) based on text found in the event explications**
  - Identified the top-most relevant terms from the event explication using inverse document frequency (IDF) on a large English language text corpus
  - Augmented the terms with associated concepts found in WordNet

# Performance of Individual Event Classifiers: EKFull



**High-level features comparable to low-level features**

# Late Fusion Models

- **No weights**
  - Arithmetic mean (AM)
  - Geometric mean (GM)
- **Fixed weights**
  - Mean average precision-weighted fusion (MAP)
  - Conditional mixture model (EM)
- **Dynamic weights**
  - Sparse conditional mixture model (SparseEM)
  - Weighted mean root
  - SVMLight
  - LibSVM
  - BBN weighting (BBN)

# Performance of Late Fusion Models: EKFull



**Simple fusion models are good enough**

# Performance of Individual Event Classifiers : EK10Ex



Fusion produces big gain in performance

# Threshold Selection Methods

- **Score@TER**
  - determined by the threshold that achieves the Target Error Ratio

- **Median score@TER**
  - for the ad hoc Ek10Ex condition only
  - median of the score@TER thresholds learned on the pre-specified events for the EK10Ex condition

- **Box-average – the average of two thresholds:**
  - the threshold that achieves P(Miss) = 50%
  - the threshold that achieves P(FA) = 4%

# SESAME MED Evaluation Runs on Progress Set

**Runs 1, 2, and 3:  Pre-specified events; EKFull; mix of extracted metadata, fusion methods, and thresholding methods**

**Run 4: Pre-specified events; EK10EX**

**Run 5: Ad hoc events; EKFull**

**Run 6: Ad hoc events; EK10Ex**

# Pre-specified Events, EKFull



**DET Plot**

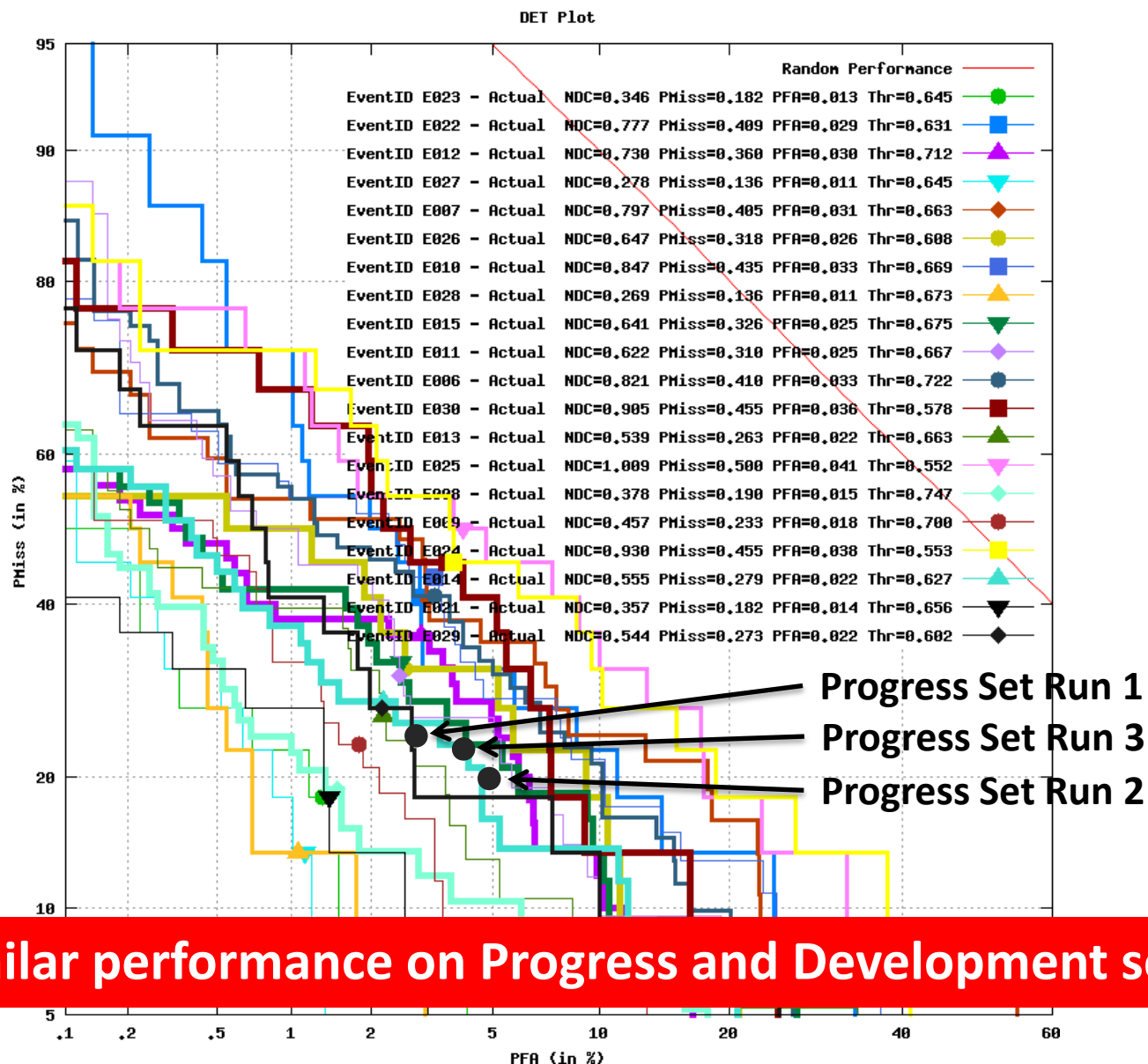| | | | | |
|---|---|---|---|---|
| | | | Random Performance | |
| EventID E023 – Actual | NDC=0.346 PMiss=0.182 PFA=0.013 Thr=0.645 | | | ● |
| EventID E022 – Actual | NDC=0.777 PMiss=0.409 PFA=0.029 Thr=0.631 | | | ■ |
| EventID E012 – Actual | NDC=0.730 PMiss=0.360 PFA=0.030 Thr=0.712 | | | ▲ |
| EventID E027 – Actual | NDC=0.278 PMiss=0.136 PFA=0.011 Thr=0.645 | | | ▼ |
| EventID E007 – Actual | NDC=0.797 PMiss=0.405 PFA=0.031 Thr=0.663 | | | ◆ |
| EventID E026 – Actual | NDC=0.647 PMiss=0.318 PFA=0.026 Thr=0.608 | | | ● |
| EventID E010 – Actual | NDC=0.847 PMiss=0.435 PFA=0.033 Thr=0.669 | | | ■ |
| EventID E028 – Actual | NDC=0.269 PMiss=0.136 PFA=0.011 Thr=0.673 | | | ▲ |
| EventID E015 – Actual | NDC=0.641 PMiss=0.326 PFA=0.025 Thr=0.675 | | | ▼ |
| EventID E011 – Actual | NDC=0.622 PMiss=0.310 PFA=0.025 Thr=0.667 | | | ◆ |
| EventID E006 – Actual | NDC=0.821 PMiss=0.410 PFA=0.033 Thr=0.722 | | | ● |
| EventID E030 – Actual | NDC=0.905 PMiss=0.455 PFA=0.036 Thr=0.578 | | | ■ |
| EventID E013 – Actual | NDC=0.539 PMiss=0.263 PFA=0.022 Thr=0.663 | | | ▲ |
| EventID E025 – Actual | NDC=1.009 PMiss=0.500 PFA=0.041 Thr=0.552 | | | ▼ |
| EventID E008 – Actual | NDC=0.378 PMiss=0.190 PFA=0.015 Thr=0.747 | | | ◆ |
| EventID E009 – Actual | NDC=0.457 PMiss=0.233 PFA=0.018 Thr=0.700 | | | ● |
| EventID E024 – Actual | NDC=0.930 PMiss=0.455 PFA=0.038 Thr=0.553 | | | ■ |
| EventID E014 – Actual | NDC=0.555 PMiss=0.279 PFA=0.022 Thr=0.627 | | | ▲ |
| EventID E021 – Actual | NDC=0.357 PMiss=0.182 PFA=0.014 Thr=0.656 | | | ▼ |
| EventID E029 – Actual | NDC=0.544 PMiss=0.273 PFA=0.022 Thr=0.602 | | | ◆ |

**Progress Set Run 1**

**Progress Set Run 3**

**Progress Set Run 2**

PMiss (in %)

PFA (in %)

**Similar performance on Progress and Development sets**

# Pre-specified Events, EK10Ex



**DET Plot**

Random Performance
EventID E023 - Actual  NDC=0.720 PMiss=0.364 PFA=0.029 Thr=0.603
EventID E022 - Actual  NDC=1.543 PMiss=0.773 PFA=0.062 Thr=0.616
EventID E012 - Actual  NDC=0.685 PMiss=0.340 PFA=0.028 Thr=0.665
EventID E027 - Actual  NDC=0.803 PMiss=0.409 PFA=0.032 Thr=0.640
EventID E007 - Actual  NDC=1.030 PMiss=0.514 PFA=0.041 Thr=0.655
EventID E026 - Actual  NDC=1.080 PMiss=0.545 PFA=0.043 Thr=0.694
EventID E013 - Actual  NDC=0.820 PMiss=0.413 PFA=0.033 Thr=0.635
EventID E020 - Actual  NDC=0.712 PMiss=0.364 PFA=0.028 Thr=0.671
EventID E015 - Actual  NDC=0.750 PMiss=0.372 PFA=0.030 Thr=0.640
EventID E011 - Actual  NDC=1.047 PMiss=0.524 PFA=0.042 Thr=0.596
EventID E006 - Actual  NDC=0.989 PMiss=0.492 PFA=0.040 Thr=0.646
EventID E030 - Actual  NDC=1.168 PMiss=0.591 PFA=0.046 Thr=0.639
EventID E013 - Actual  NDC=0.571 PMiss=0.289 PFA=0.023 Thr=0.662
EventID E025 - Actual  NDC=1.270 PMiss=0.636 PFA=0.051 Thr=0.640
EventID E008 - Actual  NDC=0.415 PMiss=0.207 PFA=0.017 Thr=0.679
EventID E009 - Actual  NDC=0.428 PMiss=0.209 PFA=0.018 Thr=0.649
EventID E024 - Actual  NDC=1.458 PMiss=0.727 PFA=0.058 Thr=0.548
EventID E014 - Actual  NDC=0.446 PMiss=0.233 PFA=0.017 Thr=0.611
EventID E021 - Actual  NDC=0.737 PMiss=0.364 PFA=0.030 Thr=0.612
EventID E029 - Actual  NDC=0.819 PMiss=0.409 PFA=0.033 Thr=0.645

**Progress Set Run 4**

**EK10Ex performance less than EKFull performance**

# Ad Hoc Events, EKFull



**Ad hoc events  as robust as pre-specified events**

# Ad Hoc Events, EK10Ex



DET Plot

Random Performance

EventID E017 - Actual  NDC=0.833 PMiss=0.432 PFA=0.032 Thr=0.663
EventID E019 - Actual  NDC=0.768 PMiss=0.381 PFA=0.031 Thr=0.610
EventID E016 - Actual  NDC=0.960 PMiss=0.465 PFA=0.040 Thr=0.631
EventID E020 - Actual  NDC=1.048 PMiss=0.581 PFA=0.037 Thr=0.553
EventID E018 - Actual  NDC=0.793 PMiss=0.400 PFA=0.031 Thr=0.649

**Progress Set Run 6**

**Threshold selection needs to improve**

PMiss (in %)

PFA (in %)

# Conclusions

- **High-level features comparable to low-level**
- **Simple average fusion good enough**
- **Similar results on Progress Set and our internal development set**
- **Ad hoc events as robust as pre-specified events**
- **Threshold selection needs to improve**

# SESAME Team for MED12

- **SRI International**
  - Gregory K. Myers, Murat Akbacak, Robert C. Bolles, J. Brian Burns, Mark Eliot, Aaron Heller, James A. Herson, Ramesh Nallapati, Stephanie Pancoast, Julien van Hout, Eric Yeh

- **University of Amsterdam**
  - Cees G.M. Snoek, Amirhossein Habibian, Dennis C. Koelma, Zhenyang Li, Masoud Mazloom, Silvia-Laura Pintea, Koen E.A. van de Sande,  Arnold W.M. Smeulders

- **University of Southern California**
  - Ram Nevatia, Sung Chun Lee, Pramod Sharma, Chen Sun, Remi Trichet

# Acknowledgement